

# 1. Business Understanding

```
In [ ]:

In [273]: #loading some relevant libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import plotly.express as px
sns.set(style='whitegrid', palette='viridis', font_scale=1.5, rc={"figure.figsize": (8, 10)})

import warnings # current version of seaborn generates a bunch of warnings that we'll ignore
warnings.filterwarnings("ignore")
```

## 2. Data Understanding and Preparation

The dataset contains medical records of 299 heart failure patients. The patients consisted of 105 women and 194 men, and their ages range between 40 and 95 years old (Table 1). All 299 patients had left ventricular systolic dysfunction and had previous heart failures that put them in classes II or IV of New York Heart Association (NYHA) classification of the stages of heart failure.

Source : <https://archive.ics.uci.edu/ml/datasets/Heart+failure+clinical+records>

```
In [89]: # Loading the dataset
hf_data = pd.read_csv("heart_failure_clinical_records_dataset.csv")

In [90]: # printing meta information
hf_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 299 entries, 0 to 298
Data columns (total 13 columns):
age                299 non-null float64
anaemia            299 non-null int64
creatinine_phosphokinase  299 non-null int64
diabetes            299 non-null int64
ejection_fraction  299 non-null int64
high_blood_pressure  299 non-null int64
platelets           299 non-null float64
serum_creatinine    299 non-null float64
serum_sodium        299 non-null int64
sex                299 non-null int64
smoking            299 non-null int64
time               299 non-null int64
DEATH_EVENT        299 non-null int64
dtypes: float64(3), int64(10)
memory usage: 30.4 KB

In [91]: # changing the datatype of age to integer
hf_data['age'] = hf_data['age'].astype('int64')
hf_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 299 entries, 0 to 298
Data columns (total 13 columns):
age                299 non-null int64
anaemia            299 non-null int64
creatinine_phosphokinase  299 non-null int64
diabetes            299 non-null int64
ejection_fraction  299 non-null int64
high_blood_pressure  299 non-null int64
platelets           299 non-null float64
serum_creatinine    299 non-null float64
serum_sodium        299 non-null int64
sex                299 non-null int64
smoking            299 non-null int64
time               299 non-null int64
DEATH_EVENT        299 non-null int64
dtypes: float64(2), int64(11)
memory usage: 30.4 KB

Quantitative Statistical Description of data

In [92]: hf_data.describe().T

Out[92]:
```

	count	mean	std	min	25%	50%	75%	max
age	299.0	60.829431	11.894997	40.0	51.0	60.0	70.0	95.0
anaemia	299.0	0.431438	0.496107	0.0	0.0	0.0	1.0	1.0
creatinine_phosphokinase	299.0	581.839465	970.287881	23.0	116.5	250.0	582.0	7861.0
diabetes	299.0	0.418950	0.494087	0.0	0.0	0.0	1.0	1.0
ejection_fraction	299.0	38.083612	11.834841	14.0	30.0	38.0	45.0	80.0
high_blood_pressure	299.0	0.521171	0.478136	0.0	0.0	0.0	1.0	1.0
platelets	299.0	263358.025264	97804.236869	25100.0	212500.0	262000.0	303500.0	850000.0
serum_creatinine	299.0	1.303880	1.034510	0.5	0.9	1.1	1.4	9.4
serum_sodium	299.0	136.625418	4.412477	113.0	134.0	137.0	140.0	148.0
sex	299.0	0.648829	0.478136	0.0	0.0	1.0	1.0	1.0
smoking	299.0	0.321070	0.467670	0.0	0.0	0.0	1.0	1.0
time	299.0	130.260870	77.614208	4.0	73.0	115.0	203.0	285.0
DEATH_EVENT	299.0	0.321070	0.467670	0.0	0.0	0.0	1.0	1.0

### Preparing a data dictionary for ease of understanding [1]

```
In [93]: labels = ['Age of the patient', 'If the patient has anaemia, 0:No 1:Yes', 'Level of the CPK enzyme in the blood',
               'If the patient has diabetes, 0:No 1:Yes', 'Percentage of blood leaving the heart at each contraction',
               'If a patient has hypertension, 0:No 1:Yes', 'Platelets in the blood',
               'Level of creatinine in the blood', 'Level of sodium in the blood', 'Woman 1: man', 'Is a smoker, 0:No, 1:Yes',
               'Follow-up period', 'If death occurred during followup period 0:No 1:Yes']

In [94]: measurement = ['Years', 'mcg/L', 'mcg/L', 'Percentage', 'Boolean', 'kiloplatelets/mL', 'mg/dL', 'm g/dL',
                       'Boolean', 'Boolean', 'Days', 'Boolean']

In [95]: hf_data_dict = pd.DataFrame({'Features':data.columns, 'Labels': labels, 'Measurement':measurement})

In [96]: hf_data_dict

Out[96]:
```

	Features	Labels	Measurement
0	age	Age of the patient	Years
1	anaemia	If the patient has anaemia, 0:No 1:Yes	Boolean
2	creatinine_phosphokinase	Level of the CPK enzyme in the blood	mcg/L
3	diabetes	If the patient has diabetes, 0:No 1:Yes	Boolean
4	ejection_fraction	Percentage of blood leaving the heart at each...	Percentage
5	high_blood_pressure	If a patient has hypertension, 0:No 1:Yes	Boolean
6	platelets	Platelets in the blood	kiloplatelets/mL
7	serum_creatinine	Level of creatinine in the blood	mg/dL
8	serum_sodium	Level of sodium in the blood	mg/dL
9	sex	0:Woman 1: man	Boolean
10	smoking	Is a smoker, 0:No, 1:Yes	Boolean
11	time	Follow-up period	Days
12	DEATH_EVENT	If death occurred during followup period 0:No 1:...	Boolean

```
In [97]: hf_data_dict.to_csv("Data Dictionary.csv", index = False) # saving for future use

Further Descriptions

Anaemia : If the hematocrit levels were less than 36% then classified as No. Hematocrit is the volume percentage of Red Blood Cells in blood[2].

Creatinine Phosphokinase: Acronymed as CPK, this enzyme is an indicator of heart failure.CPK flows into the blood upon damage to a muscle tissue. Higher levels of CPK in blood might indicate a heart failure[3].

Ejection Fraction: Ejection fraction, expressed as a percentage measures how much blood is pumped out by the left ventricle with each contraction. This indicates how well the heart is at pumping out blood. It aids in diagnosis and track heart failure[4].

Platelets : Platelets are specialized disk-shaped cells in the blood stream that are involved in the formation of blood clots that play an important role in heart attacks, strokes, and peripheral vascular disease.Platelets can detect a disruption in the lining of a blood vessel and react to build a wall to stop bleeding[5].

Serum Creatinine: Serum Creatinine is a blood chemical waste product generated from muscle breakdown. It is associated with functioning of kidney. High levels of serum creatinine indicate renal dysfunction.Renal dysfunction has shown to influence heart failure[6].

Serum Sodium: Sodium is a type of electrolyte. An electrolyte is a charged mineral that aids correct functioning of muscles and nerves. This is done by maintaining a balance in fluid levels and chemicals(acid & bases) in the body. An abnormally low level of sodium might be caused by heart failure [7].

Separating the features into numerical and categorical

The data has been separated according to the nature of features. This will help in better analysis. Which type of plots to use for visualising determined by the type of data. Also, the type of hypothesis tests to use could be easily identified once there is a clarity regarding the data type.

According to our dataset:

Numerical Attributes - age, creatinine_phosphokinase, ejection_fraction,platelets, serum_creatinine, serum_sodium,time

Categorical Attributes - anaemia, diabetes, high_blood_pressure, sex, smoking, DEATH_EVENT
```

```
In [98]: numerical = ['age', 'creatinine_phosphokinase', 'ejection_fraction', 'platelets', 'serum_creatinine', 's
         serum_sodium', 'time']
categorical = ['anaemia','diabetes', 'high_blood_pressure', 'sex', 'smoking', 'DEATH_EVENT']

In [99]: hf_data_EDA = hf_data[numerical + categorical]
hf_data_EDA.head()

Out[99]:
```

	age	creatinine_phosphokinase	ejection_fraction	platelets	serum_creatinine	serum_sodium	time	anaemia	diabetes	high_blood_prs
0	75	582	20	265000.0	1.9	130	4	0	0	0
1	55	7861	38	263356.03	1.1	136	6	0	0	0
2	65	146	20	162000.0	1.3	129	7	0	0	0
3	50	111	20	210000.0	1.9	137	7	1	0	0
4	65	160	20	327000.0	2.7	116	8	1	1	0

## 2.1 Exploratory Data Analysis

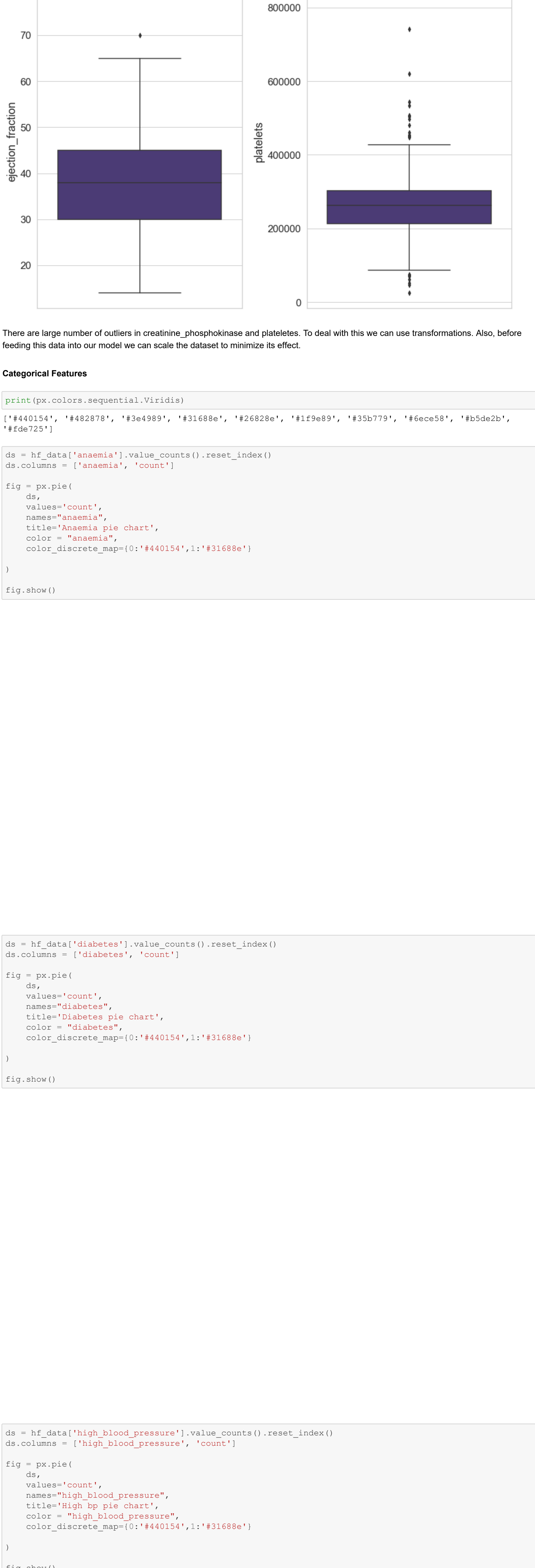
```
In [ ]:

2.1.1 Univariate Analysis

Numerical Features
```

```
In [120]: hf_data[numerical].hist(bins=10,figsize=(20, 28), layout=(4, 2))

Out[120]: array([[<matplotlib.axes._subplots.AxesSubplot object at 0x000001AC66FA3710>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x000001AC86FD9A20>],
                [<matplotlib.axes._subplots.AxesSubplot object at 0x000001AC87007400>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x000001AC870370AD>],
                [<matplotlib.axes._subplots.AxesSubplot object at 0x000001AC87066780>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x000001AC870AB160>],
                [<matplotlib.axes._subplots.AxesSubplot object at 0x000001AC8674C438>,
                  <matplotlib.axes._subplots.AxesSubplot object at 0x000001AC670B07D0>]],
              dtype=object)
```



```
In [134]: fig, ax = plt.subplots(2, 2, figsize=(12, 16))
for variable, subplot in zip(numerical, ax.flatten()):
    sns.boxplot(y=hf_data_EDA[variable], ax=subplot)
fig.tight_layout()
```



There are large number of outliers in creatinine\_phosphokinase and platelets. To deal with this we can use transformations. Also, before feeding this data into our model we can scale the dataset to minimize its effect.

### Categorical Features

```
In [238]: print(px.colors.sequential.Viridis)

['#440154', '#482878', '#3e4989', '#31688e', '#26828e', '#1f9e89', '#35b779', '#6ece58', '#b5de2b',
 '#d6e725']

In [240]: ds = hf_data['anaemia'].value_counts().reset_index()
ds.columns = ['anaemia', 'count']

fig = px.pie(
    ds,
    values='count',
    names='anaemia',
    title='Anaemia pie chart',
    color = "anaemia",
    color_discrete_map={0:'#440154',1:'#31688e'}
)

fig.show()
```

```
In [242]: ds = hf_data['diabetes'].value_counts().reset_index()
ds.columns = ['diabetes', 'count']

fig = px.pie(
    ds,
    values='count',
    names='diabetes',
    title='Diabetes pie chart',
    color = "diabetes",
    color_discrete_map={0:'#440154',1:'#31688e'}
)

fig.show()
```

```
In [244]: ds = hf_data['high_blood_pressure'].value_counts().reset_index()
ds.columns = ['high_blood_pressure', 'count']

fig = px.pie(
    ds,
    values='count',
    names='high blood pressure',
    title='High bp pie chart',
    color = "high_blood_pressure",
    color_discrete_map={0:'#440154',1:'#31688e'}
)

fig.show()
```

```
In [245]: ds = hf_data['sex'].value_counts().reset_index()
ds.columns = ['sex', 'count']

fig = px.pie(
    ds,
    values='count',
    names='sex',
    title='Sex pie chart',
    color = "sex",
    color_discrete_map={0:'#440154',1:'#31688e'}
)

fig.show()
```

```
In [246]: ds = hf_data['smoking'].value_counts().reset_index()
ds.columns = ['smoking', 'count']

fig = px.pie(
    ds,
    values='count',
    names='smoking',
    title='Smoking pie chart',
    color = "smoking",
    color_discrete_map={0:'#440154',1:'#31688e'}
)

fig.show()
```

```
In [248]: ds = hf_data['DEATH_EVENT'].value_counts().reset_index()
ds.columns = ['DEATH_EVENT', 'count']

fig = px.pie(
    ds,
    values='count',
    names='DEATH_EVENT',
    title='Death Event pie chart',
    color = "DEATH_EVENT",
    color_discrete_map={0:'#440154',1:'#31688e'}
)

fig.show()
```

```
In [249]: #fig, ax = plt.subplots(2, 2, figsize=(20, 18))
#sns.heatmap(hf_data_EDA[numerical], corr=True, annot=True)
#sns.countplot(hf_data_EDA[variable], ax=subplot)
```

## 2.1.2 Bivariate Analysis

### Numerical Features

#### Correlation

```
In [247]: hf_data_EDA[numerical].corr()

Out[192]:
```

	age	creatinine_phosphokinase	ejection_fraction	platelets	serum_creatinine	serum_sodium	time
age	1.000000	-0.081406	-0.081406	0.060195	-0.052475	0.159237	-0.045912
creatinine_phosphokinase	-0.081406	1.000000	-0.044080	0.024463	-0.016408	0.059550	-0.009346
ejection_fraction	-0.081406	-0.044080	1.000000	0.072177	-0.011302	0.175902	0.041729
platelets	-0.052475	0.024463	0.072177	1.000000	-0.041198	0.062125	0.010514
serum_creatinine	0.159237	-0.016408	-0.011302	-0.041198	1.000000	-0.189095	-0.149315
serum_sodium	-0.045912	0.059550	0.175902	0.062125	-0.189095	1.000000	0.087640
time	-0.224265	-0.009346	0.041729	0.010514	-0.149315	0.087640	1.000000

```
In [264]: plt.subplots(figsize=(16,8))
sns.heatmap(hf_data_EDA[numerical].corr(), cmap = "viridis", annot=True)

Out[264]: <matplotlib.axes._subplots.AxesSubplot at 0x1ac94951da0>
```

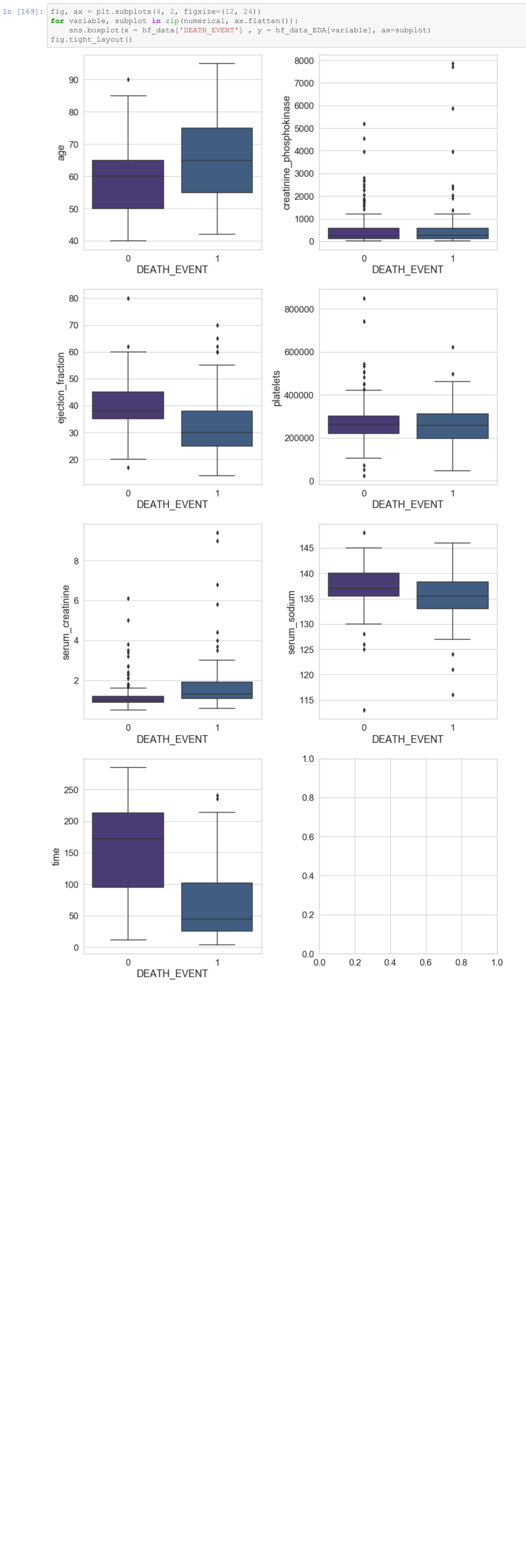


The numeric Independent variables are not highly correlated. This is good as it indicates absence of multicollinearity

### Analysis with target variable

Our objective is to predict Death Event. Analysis of the features with target variable could give us insights into the effect they could have on our predictions.



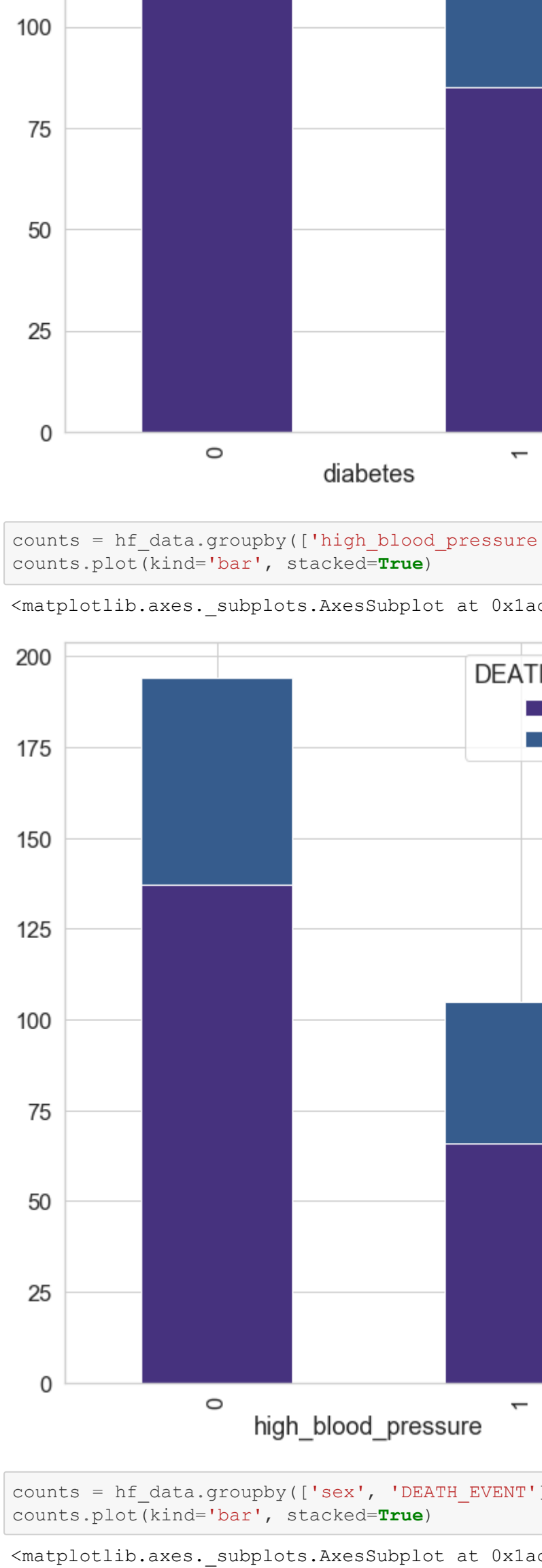


## Categorical Features

```
In [169]: hf_data_EDA[categorical].columns
Out[169]: Index(['anaemia', 'diabetes', 'high_blood_pressure', 'sex', 'smoking',
              'DEATH_EVENT'],
              dtype='object')
```

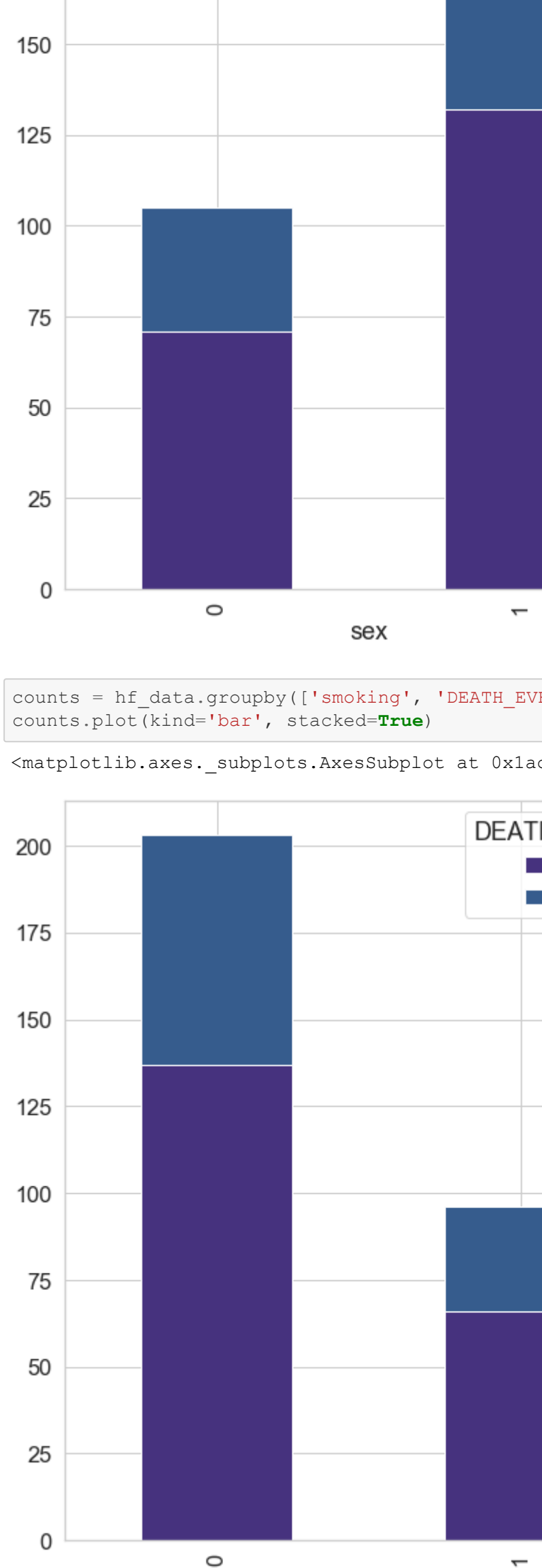
```
In [172]: counts = hf_data.groupby(['anaemia', 'DEATH_EVENT']).anaemia.count().unstack()
counts.plot(kind='bar', stacked=True)
```

```
Out[172]: <matplotlib.axes._subplots.AxesSubplot at 0x1ac858dc18>
```



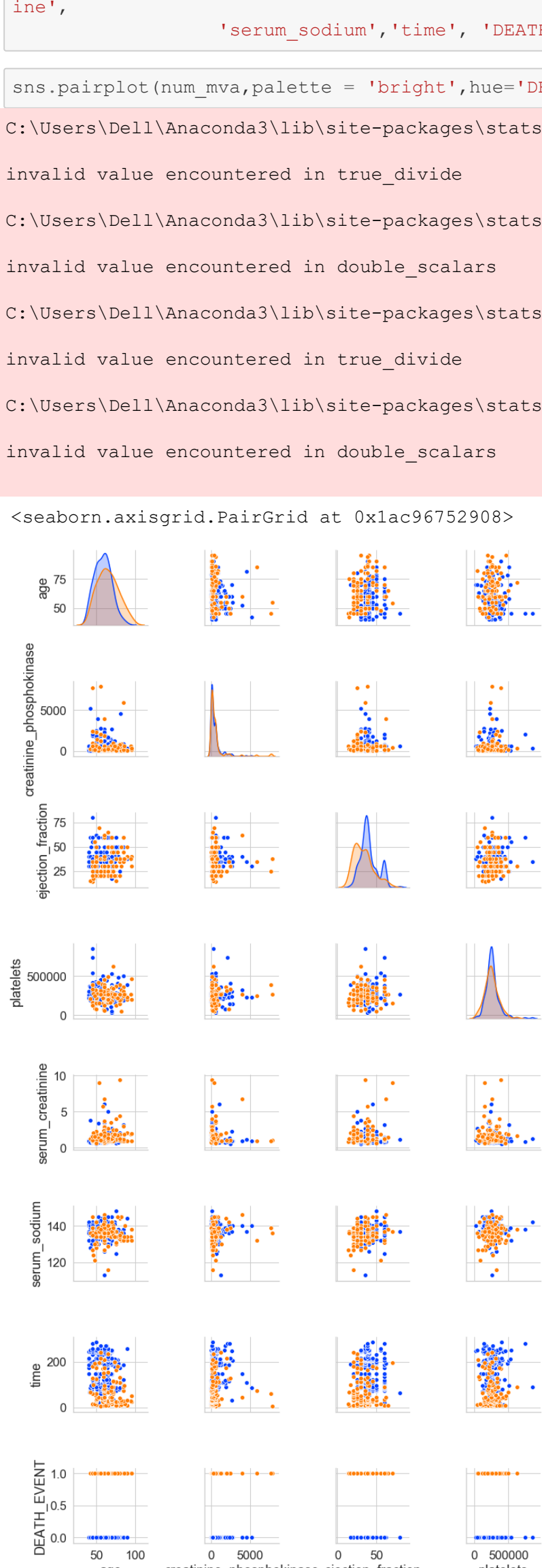
```
In [171]: counts = hf_data.groupby(['smoking', 'DEATH_EVENT']).diabetes.count().unstack()
counts.plot(kind='bar', stacked=True)
```

```
Out[171]: <matplotlib.axes._subplots.AxesSubplot at 0x1ac8510ac8>
```



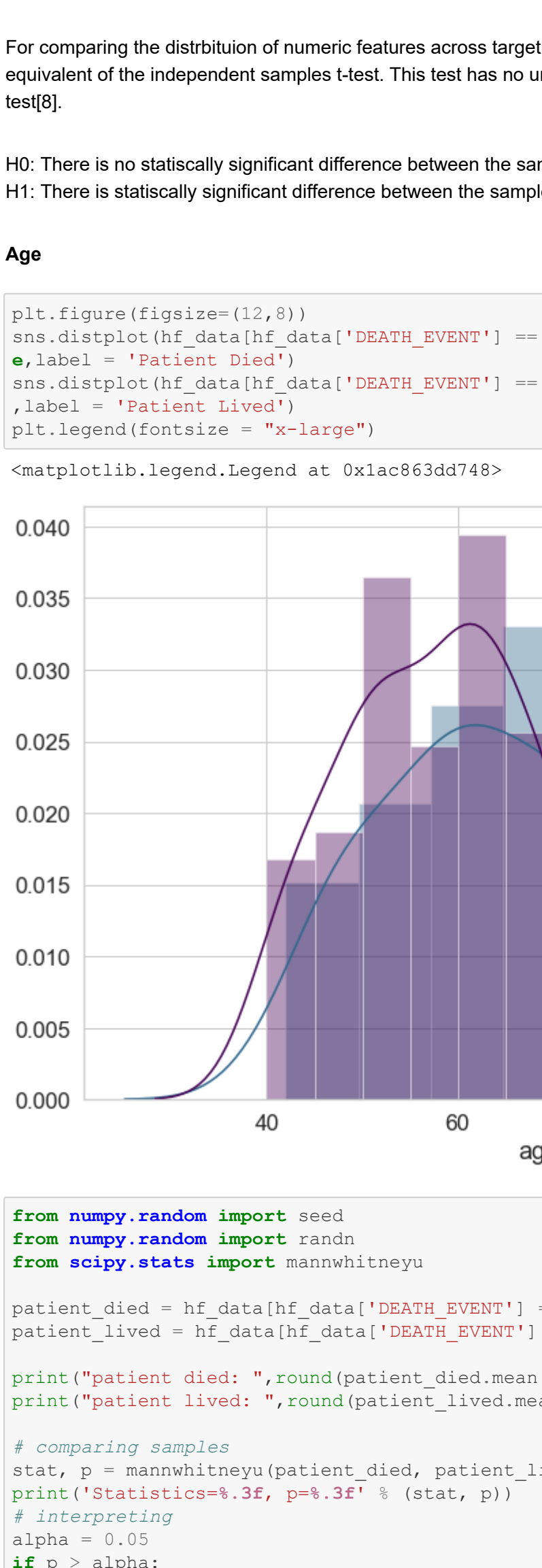
```
In [173]: counts = hf_data.groupby(['high_blood_pressure', 'DEATH_EVENT']).high_blood_pressure.count().unstack()
counts.plot(kind='bar', stacked=True)
```

```
Out[173]: <matplotlib.axes._subplots.AxesSubplot at 0x1ac8b8f8ba>
```



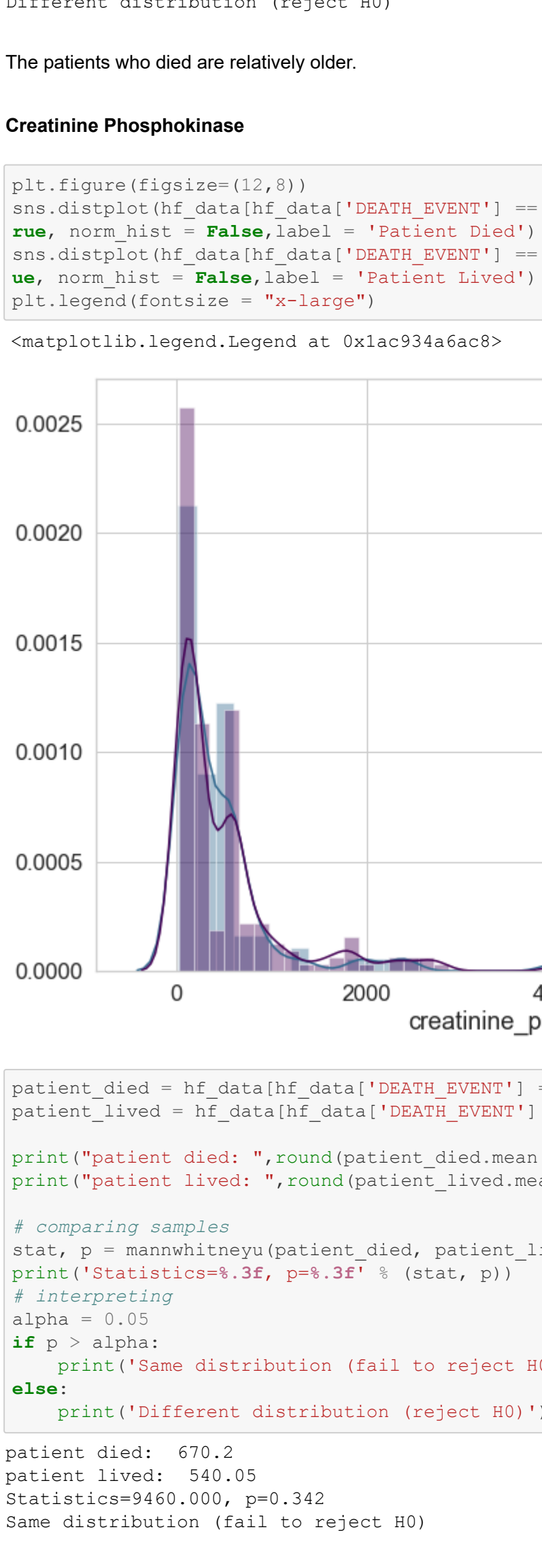
```
In [174]: counts = hf_data.groupby(['sex', 'DEATH_EVENT']).sex.count().unstack()
counts.plot(kind='bar', stacked=True)
```

```
Out[174]: <matplotlib.axes._subplots.AxesSubplot at 0x1ac8bfe668>
```



```
In [175]: counts = hf_data.groupby(['smoking', 'DEATH_EVENT']).smoking.count().unstack()
counts.plot(kind='bar', stacked=True)
```

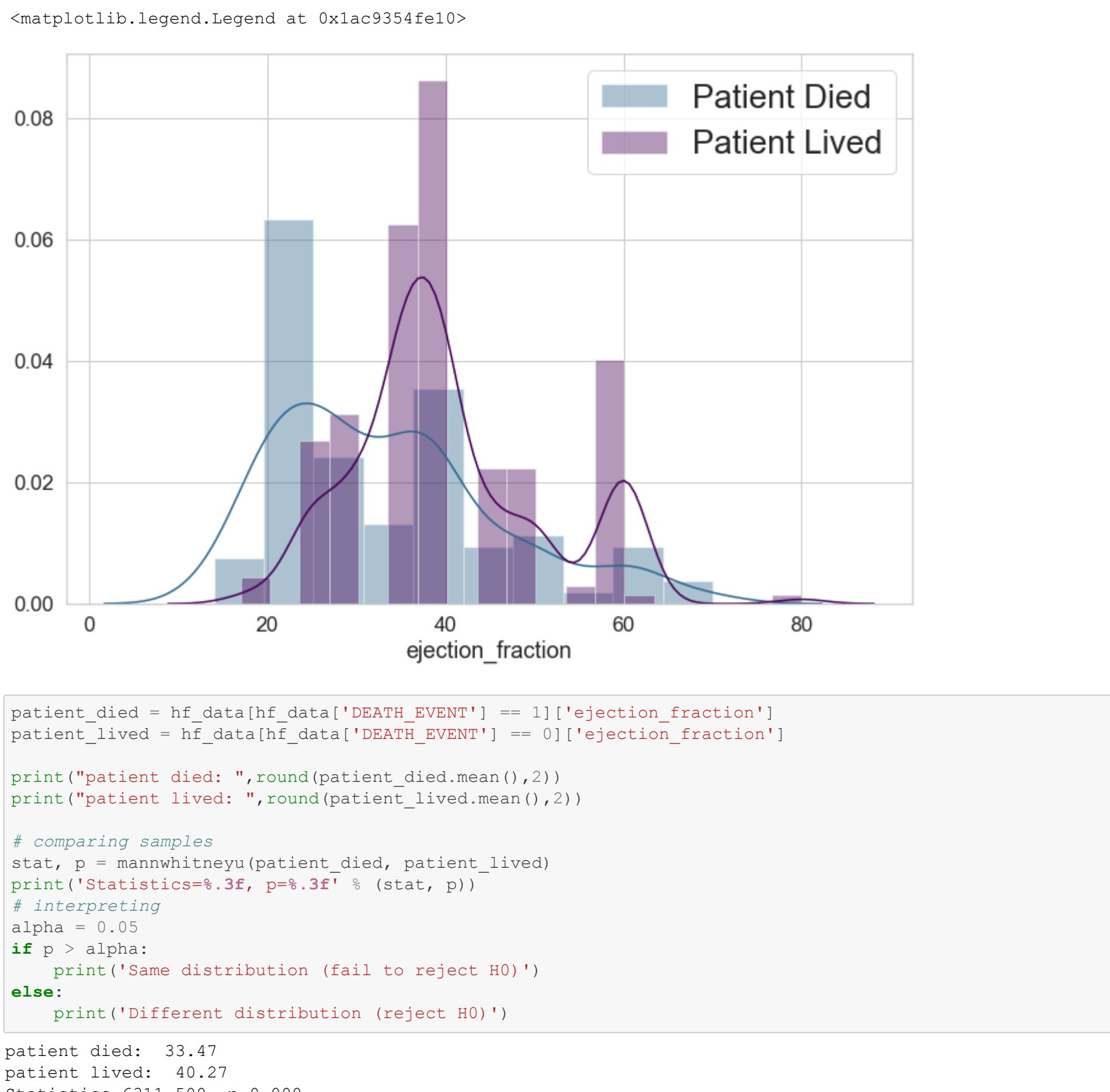
```
Out[175]: <matplotlib.axes._subplots.AxesSubplot at 0x1ac9a5c5828>
```



## 2.1.3 Multivariate Analysis

```
In [256]: num_mva = hf_data[['age', 'creatinine_phosphokinase', 'ejection_fraction', 'platelets', 'serum_creatinine',
                          'serum_sodium', 'time', 'DEATH_EVENT']]
```

```
In [257]: sns.pairplot(num_mva, palette = 'bright', hue='DEATH_EVENT')
```



## 2.2 Hypothesis testing

### 2.2.1 Numerical Variables

For comparing the distribution of numeric features across target variable we have used the Mann-Whitney test. It is the non-parametric equivalent of the independent samples t-test. This test has no underlying assumption regarding the normal distribution of data unlike t-test.

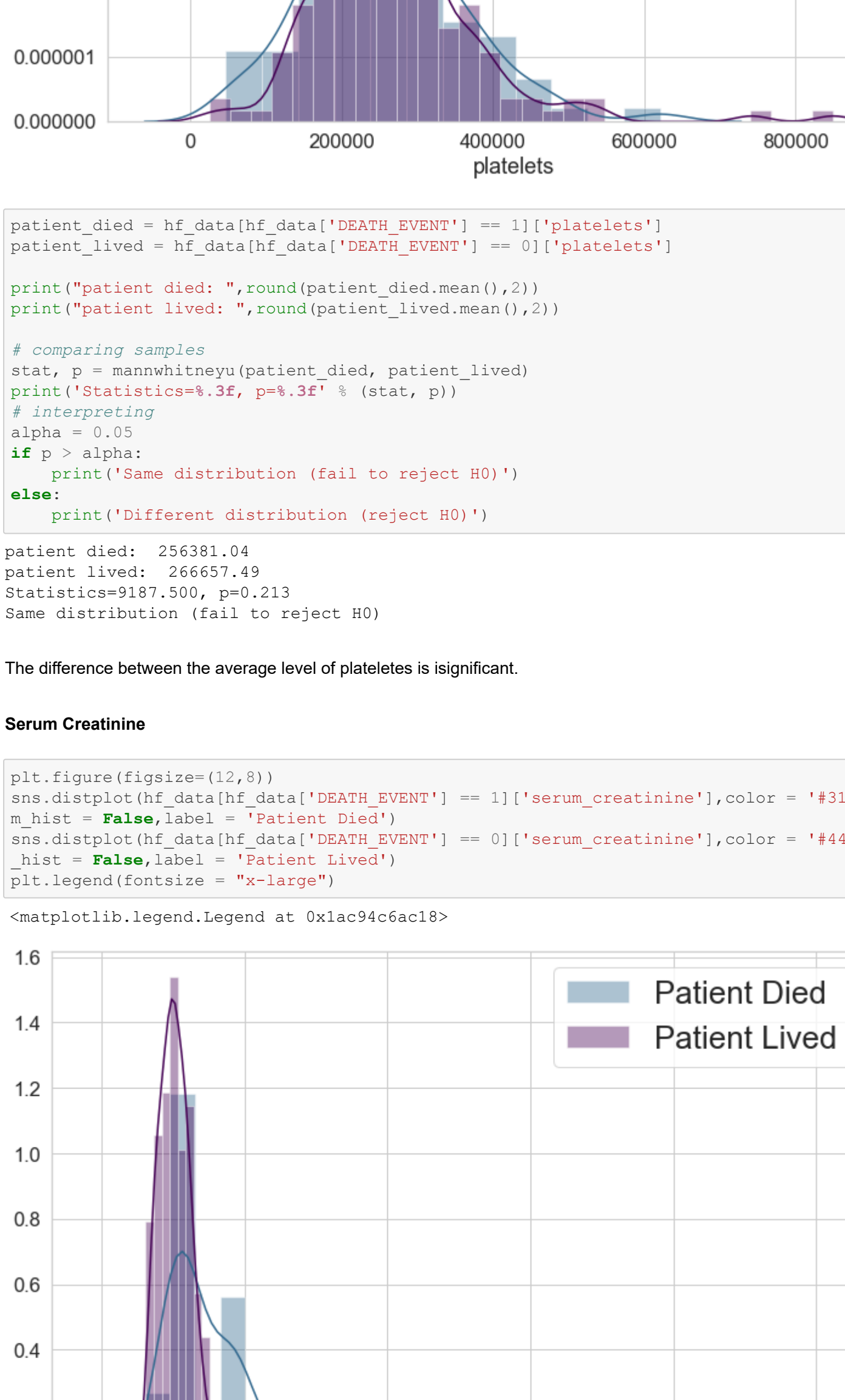
H0: There is no statistically significant difference between the samples

H1: There is statistically significant difference between the samples

#### Age

```
In [272]: plt.figure(figsize=(12,8))
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 1]['age'], color = '#31688e', kde=True, norm_hist = False)
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 0]['age'], color = '#440154', kde=True, norm_hist = False)
plt.legend(fontsize = 'x-large')
```

```
Out[272]: <matplotlib.legend.Legend at 0x1ac936dd748>
```



```
In [296]: from numpy.random import seed
from numpy.random import randn
from scipy.stats import mannwhitneyu

patient_died = hf_data[hf_data['DEATH_EVENT'] == 1]['age']
patient_lived = hf_data[hf_data['DEATH_EVENT'] == 0]['age']

print("patient died: ", round(patient_died.mean(), 2))
print("patient lived: ", round(patient_lived.mean(), 2))

# comparing samples
stat, p = mannwhitneyu(patient_died, patient_lived)
print("Statistics=%.3f, p=%.3f" % (stat, p))
# interpreting
alpha = 0.05
if p > alpha:
    print("Same distribution (fail to reject H0)")
else:
    print("Different distribution (reject H0)")
```

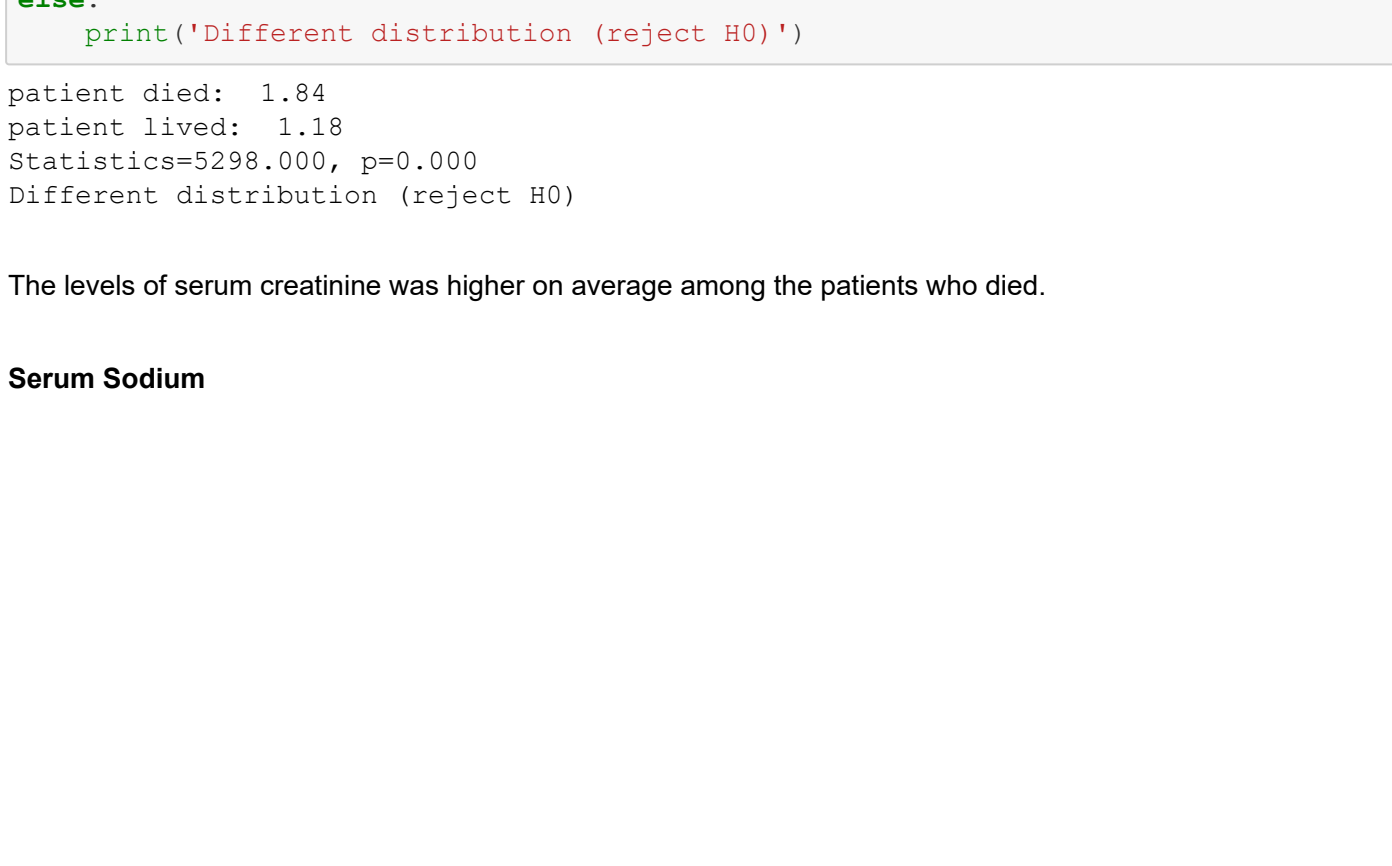
patient died: 65.21  
patient lived: 58.76  
Statistics=7124.500, p=0.000  
Different distribution (reject H0)

The patients who died are relatively older.

#### Creatinine Phosphokinase

```
In [277]: plt.figure(figsize=(12,8))
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 1]['creatinine_phosphokinase'], color = '#31688e', kde=True, norm_hist = False)
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 0]['creatinine_phosphokinase'], color = '#440154', kde=True, norm_hist = False)
plt.legend(fontsize = 'x-large')
```

```
Out[277]: <matplotlib.legend.Legend at 0x1ac9346ac8>
```



```
In [297]: patient_died = hf_data[hf_data['DEATH_EVENT'] == 1]['creatinine_phosphokinase']
patient_lived = hf_data[hf_data['DEATH_EVENT'] == 0]['creatinine_phosphokinase']

print("patient died: ", round(patient_died.mean(), 2))
print("patient lived: ", round(patient_lived.mean(), 2))

# comparing samples
stat, p = mannwhitneyu(patient_died, patient_lived)
print("Statistics=%.3f, p=%.3f" % (stat, p))
# interpreting
alpha = 0.05
if p > alpha:
    print("Same distribution (fail to reject H0)")
else:
    print("Different distribution (reject H0)")
```

patient died: 670.2  
patient lived: 540.05  
Statistics=9460.000, p=0.342  
Same distribution (fail to reject H0)

The level of CPK enzyme is not higher among patients who died.

#### Ejection Fraction

```
In [280]: plt.figure(figsize=(12,8))
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 1]['ejection_fraction'], color = '#31688e', kde=True, norm_hist = False)
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 0]['ejection_fraction'], color = '#440154', kde=True, norm_hist = False)
plt.legend(fontsize = 'x-large')
```

```
Out[280]: <matplotlib.legend.Legend at 0x1ac9354fe0>
```



```
In [299]: patient_died = hf_data[hf_data['DEATH_EVENT'] == 1]['ejection_fraction']
patient_lived = hf_data[hf_data['DEATH_EVENT'] == 0]['ejection_fraction']

print("patient died: ", round(patient_died.mean(), 2))
print("patient lived: ", round(patient_lived.mean(), 2))

# comparing samples
stat, p = mannwhitneyu(patient_died, patient_lived)
print("Statistics=%.3f, p=%.3f" % (stat, p))
# interpreting
alpha = 0.05
if p > alpha:
    print("Same distribution (fail to reject H0)")
else:
    print("Different distribution (reject H0)")
```

patient died: 33.47  
patient lived: 40.27  
Statistics=8311.500, p=0.000  
Different distribution (reject H0)

The ejection fraction was lower among the patients who died.

#### Platelets

```
In [301]: plt.figure(figsize=(12,8))
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 1]['platelets'], color = '#31688e', kde=True, norm_hist = False)
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 0]['platelets'], color = '#440154', kde=True, norm_hist = False)
plt.legend(fontsize = 'x-large')
```

```
Out[301]: <matplotlib.legend.Legend at 0x1ac934fcd80>
```



```
In [305]: patient_died = hf_data[hf_data['DEATH_EVENT'] == 1]['platelets']
patient_lived = hf_data[hf_data['DEATH_EVENT'] == 0]['platelets']

print("patient died: ", round(patient_died.mean(), 2))
print("patient lived: ", round(patient_lived.mean(), 2))

# comparing samples
stat, p = mannwhitneyu(patient_died, patient_lived)
print("Statistics=%.3f, p=%.3f" % (stat, p))
# interpreting
alpha = 0.05
if p > alpha:
    print("Same distribution (fail to reject H0)")
else:
    print("Different distribution (reject H0)")
```

patient died: 256381.04  
patient lived: 266657.49  
Statistics=9187.500, p=0.213  
Same distribution (fail to reject H0)

The difference between the average level of platelets is insignificant.

#### Serum Creatinine

```
In [303]: plt.figure(figsize=(12,8))
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 1]['serum_creatinine'], color = '#31688e', kde=True, norm_hist = False)
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 0]['serum_creatinine'], color = '#440154', kde=True, norm_hist = False)
plt.legend(fontsize = 'x-large')
```

```
Out[303]: <matplotlib.legend.Legend at 0x1ac94c6ac18>
```



```
In [307]: patient_died = hf_data[hf_data['DEATH_EVENT'] == 1]['serum_creatinine']
patient_lived = hf_data[hf_data['DEATH_EVENT'] == 0]['serum_creatinine']

print("patient died: ", round(patient_died.mean(), 2))
print("patient lived: ", round(patient_lived.mean(), 2))

# comparing samples
stat, p = mannwhitneyu(patient_died, patient_lived)
print("Statistics=%.3f, p=%.3f" % (stat, p))
# interpreting
alpha = 0.05
if p > alpha:
    print("Same distribution (fail to reject H0)")
else:
    print("Different distribution (reject H0)")
```

patient died: 1.94  
patient lived: 1.18  
Statistics=5298.000, p=0.000  
Different distribution (reject H0)

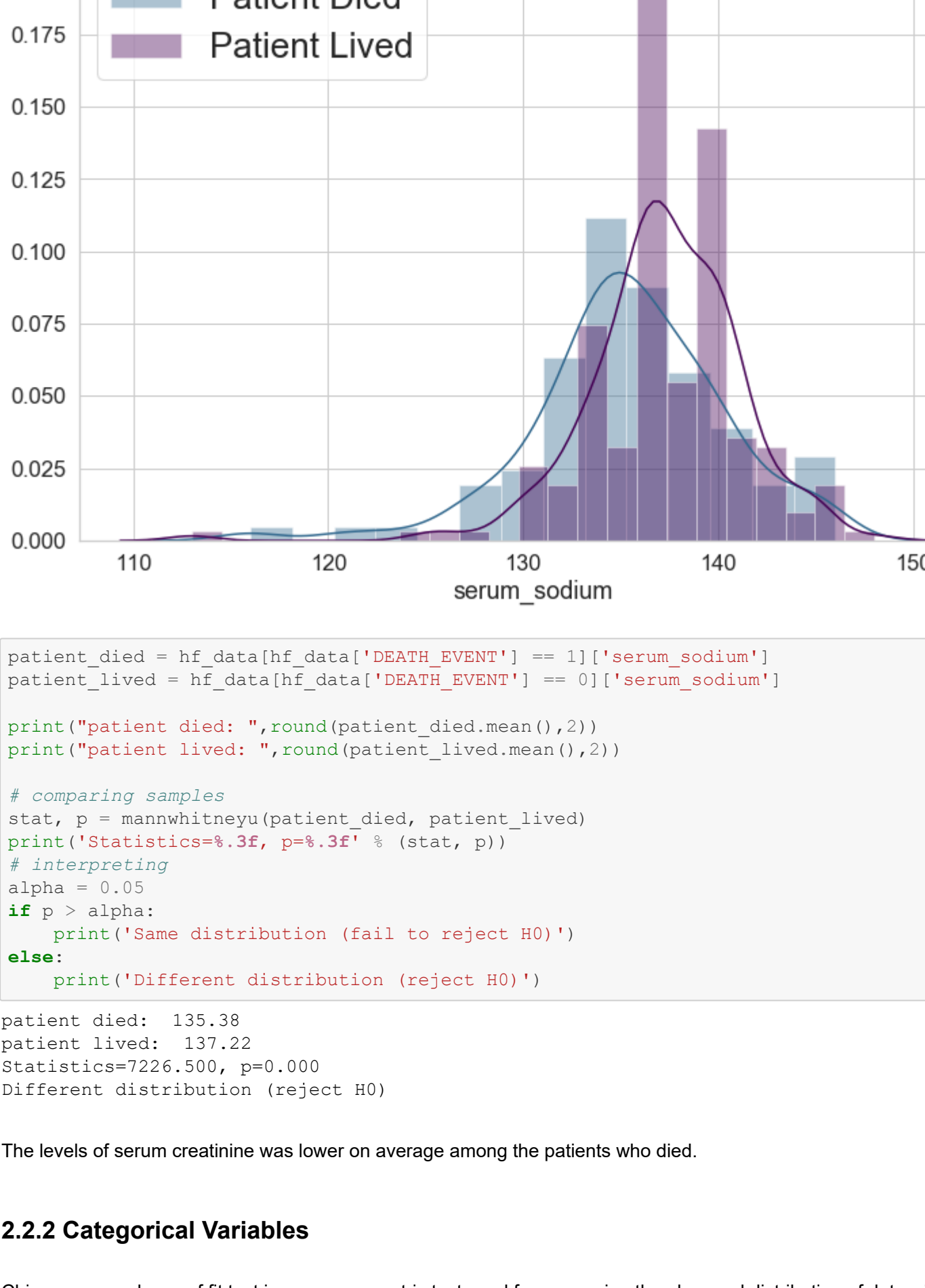
The levels of serum creatinine was higher on average among the patients who died.

#### Serum Sodium



```
In [335]: plt.figure(figsize=(12,8))
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 1]['serum_sodium'],color = '#31688e', kde=True, norm_his
st = False, label = 'Patient Died')
sns.distplot(hf_data[hf_data['DEATH_EVENT'] == 0]['serum_sodium'],color = '#440154',kde=True, norm_his
t = False, label = 'Patient Lived')
plt.legend(fontsize = "x-large")
```

Out[335]: <matplotlib.legend.Legend at 0x1ac997836d8>



```
In [309]: patient_died = hf_data[hf_data['DEATH_EVENT'] == 1]['serum_sodium']
patient_lived = hf_data[hf_data['DEATH_EVENT'] == 0]['serum_sodium']
print("patient died: ",round(patient_died.mean(),2))
print("patient lived: ",round(patient_lived.mean(),2))

# comparing samples
stat, p = mannwhitneyu(patient_died, patient_lived)
print("Statistics=%3f, p=%3f" % (stat, p))
# interpreting
alpha = 0.05
if p > alpha:
    print("Same distribution (fail to reject H0)")
else:
    print("Different distribution (reject H0)")

patient died: 135.38
patient lived: 137.22
Statistics=225.500, p=0.000
Different distribution (reject H0)
```

The levels of serum creatinine was lower on average among the patients who died.

## 2.2.2 Categorical Variables

Chi-square goodness of fit test is a non-parametric test used for comparing the observed distribution of data with the expected distribution of the data to decide whether there is any statistically significant difference between the observed distribution and a theoretical distribution.

H0: There is no difference between observed and expected frequencies

HA: There is difference between observed and expected frequencies

### Anaemia

```
In [313]: tab = pd.crosstab(index = hf_data['anaemia'], columns = hf_data['DEATH_EVENT'])
print(tab)

from scipy.stats import chi2_contingency
stat, pvalue, dof, expected = chi2_contingency(tab, correction = True)
print("chi-sq-statistic: ", round(stat, 4), " p-value: ", round(pvalue, 4), " deg of freedom: ", dof )

if pvalue < 0.05:
    print ('Reject Null Hypothesis')
else:
    print ('Retain Null Hypothesis')
```

DEATH\_EVENT 0 1  
anaemia 0 120 50  
1 13 46  
chi-sq-statistic: 1.0422 p-value: 0.3073 deg of freedom: 1  
Retain Null Hypothesis

The relationship between anaemia and patient death is insignificant.

### Diabetes

```
In [315]: tab = pd.crosstab(index = hf_data['diabetes'], columns = hf_data['DEATH_EVENT'])
print(tab)

from scipy.stats import chi2_contingency
stat, pvalue, dof, expected = chi2_contingency(tab, correction = True)
print("chi-sq-statistic: ", round(stat, 4), " p-value: ", round(pvalue, 4), " deg of freedom: ", dof )

if pvalue < 0.05:
    print ('Reject Null Hypothesis')
else:
    print ('Retain Null Hypothesis')
```

DEATH\_EVENT 0 1  
diabetes 0 118 56  
1 45 40  
chi-sq-statistic: 0.0085 p-value: 0.9267 deg of freedom: 1  
Retain Null Hypothesis

The relationship between diabetes and patient death is insignificant.

### High Blood Pressure

```
In [317]: tab = pd.crosstab(index = hf_data['high_blood_pressure'], columns = hf_data['DEATH_EVENT'])
print(tab)

from scipy.stats import chi2_contingency
stat, pvalue, dof, expected = chi2_contingency(tab, correction = True)
print("chi-sq-statistic: ", round(stat, 4), " p-value: ", round(pvalue, 4), " deg of freedom: ", dof )

if pvalue < 0.05:
    print ('Reject Null Hypothesis')
else:
    print ('Retain Null Hypothesis')
```

DEATH\_EVENT 0 1  
high\_blood\_pressure 0 137 57  
1 66 39  
chi-sq-statistic: 1.5435 p-value: 0.2141 deg of freedom: 1  
Retain Null Hypothesis

The relationship between high blood pressure and patient death is insignificant.

### Sex

```
In [319]: tab = pd.crosstab(index = hf_data['sex'], columns = hf_data['DEATH_EVENT'])
print(tab)

from scipy.stats import chi2_contingency
stat, pvalue, dof, expected = chi2_contingency(tab, correction = True)
print("chi-sq-statistic: ", round(stat, 4), " p-value: ", round(pvalue, 4), " deg of freedom: ", dof )

if pvalue < 0.05:
    print ('Reject Null Hypothesis')
else:
    print ('Retain Null Hypothesis')
```

DEATH\_EVENT 0 1  
sex 0 71 34  
1 132 62  
chi-sq-statistic: 0.0003 p-value: 0.9561 deg of freedom: 1  
Retain Null Hypothesis

The relationship between sex and patient death is insignificant.

### Smoking

```
In [321]: tab = pd.crosstab(index = hf_data['smoking'], columns = hf_data['DEATH_EVENT'])
print(tab)

from scipy.stats import chi2_contingency
stat, pvalue, dof, expected = chi2_contingency(tab, correction = True)
print("chi-sq-statistic: ", round(stat, 4), " p-value: ", round(pvalue, 4), " deg of freedom: ", dof )

if pvalue < 0.05:
    print ('Reject Null Hypothesis')
else:
    print ('Retain Null Hypothesis')
```

DEATH\_EVENT 0 1  
smoking 0 137 66  
1 66 30  
chi-sq-statistic: 0.0073 p-value: 0.9318 deg of freedom: 1  
Retain Null Hypothesis

The relationship between smoking and patient death is insignificant.

## 3. Modeling & Evaluation

The aim of our model would be to predict the event of patient's death. Since the outcome is categorical in nature, classification algorithms have been applied.

### Separating Dependent and Independent variables

```
In [336]: # Followup period would automatically be less for patients who have died so it is redundant
# Death Event is the Dependent Variable
features = hf_data.drop(['DEATH_EVENT', 'time'],axis=1)
X_features = list(features.columns)

Y = hf_data.DEATH_EVENT
```

### 3.1 Logistic Regression

Logistic Regression using statsmodel package has been used. The statsmodel package includes well tabulated summary for the model.

Advantages:

1. Easy interpretation
2. Probabilities of outcome can be estimated
3. Robust from overfitting
4. No assumptions pertaining to the distribution

```
In [336]: # Adding a constant for Logistic model to work while using statsmodels package
import statsmodels.api as sm
X = sm.add_constant(hf_data[X_features])
X.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 239 entries, 0 to 238
Data columns (total 12 columns):
const                239 non-null float64
age                  239 non-null int64
anaemia              239 non-null int64
creatinine_phosphokinase  239 non-null int64
diabetes             239 non-null int64
ejection_fraction    239 non-null int64
high_blood_pressure  239 non-null int64
platelets            239 non-null float64
serum_creatinine      239 non-null float64
serum_sodium         239 non-null int64
sex                  239 non-null int64
smoking              239 non-null int64
dtype: float64(3), int64(9)
memory usage: 28.1 KB
```

```
In [434]: # Splitting into Train and Validation Sets
# 80% Train, 20% Validation
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X,Y,test_size = 0.20,random_state = 42)
```

```
In [435]: # Building the logistic regression model
import statsmodels.api as sm
logit = sm.Logit(y_train, X_train)
logit_model = logit.fit()

Optimization terminated successfully.
Current function value: 0.455891
Iterations: 7
```

```
In [436]: # Printing Model Summary
logit_model.summary2()
```

Out[436]:

	Model:	Logit	Pseudo R-squared:	0.261
Dependent Variable:	DEATH_EVENT	AIC:	241.9110	
	Date: 2020-10-11 03:43	BIC:	283.6336	
No. Observations:	239	Log-Likelihood:	-108.96	
DF Model:	11	LL-Nucll:	-145.40	
DF Residuals:	227	LLR p-value:	3.4486e-11	
Converged:	1.0000	Scale:	1.0000	
No. iterations:	7.0000			
	Coeff.	Std. Err.	z	P> z  [0.025 0.975]
const	2.8202	5.3507	0.5271	0.5981 -7.6669 13.3073
age	0.0681	0.0156	4.3669	0.0000 0.0375 0.0987
anaemia	0.3699	0.2517	1.0517	0.2930 -0.3195 1.0593
creatinine_phosphokinase	0.0002	0.0002	1.3940	0.1633 -0.0001 0.0006
diabetes	0.2671	0.3548	0.8372	0.4025 -0.3894 0.9925
ejection_fraction	-0.0784	0.0182	-4.3058	0.0000 -0.1141 -0.0427
high_blood_pressure	0.4646	0.3602	1.2901	0.1970 -0.2413 1.1705
platelets	-0.0000	0.0000	-0.8947	0.3710 -0.0000 0.0000
serum_creatinine	0.7502	0.1955	3.8366	0.0001 0.3670 1.1335
serum_sodium	-0.0440	0.0383	-1.1485	0.2508 -0.1192 0.0311
sex	-0.6644	0.4042	-1.6438	0.1002 -1.4566 1.1278
smoking	0.1689	0.4128	0.4091	0.6824 -0.6402 0.9780

```
In [437]: # Model Diagnostics
def get_significant_vars (lm):
    var_p_vals_df = pd.DataFrame( { "actual": y_test,
    var_p_vals_df["vars"] = var_p_vals_df.index
    var_p_vals_df.columns = ['pvals', 'vars']
    return list( var_p_vals_df[var_p_vals_df.pvals <= 0.05]['vars'] )
```

```
In [438]: # Significant Variables
significant_vars = get_significant_vars( logit_model )
significant_vars
```

Out[438]: ['age', 'ejection\_fraction', 'serum\_creatinine']

```
In [439]: # Rebuilding the model with significant variables
final_logit = sm.Logit( y_train, sm.add_constant( X_train[significant_vars] ) ).fit()

Optimization terminated successfully.
Current function value: 0.475243
Iterations: 6
```

```
In [440]: # Printing Model Summary
final_logit.summary2()
```

Out[440]:

	Model:	Logit	Pseudo R-squared:	0.219
Dependent Variable:	DEATH_EVENT	AIC:	235.1710	
	Date: 2020-10-11 03:43	BIC:	249.0768	
No. Observations:	239	Log-Likelihood:	-113.59	
DF Model:	3	LL-Nucll:	-145.40	
DF Residuals:	235	LLR p-value:	9.8669e-14	
Converged:	1.0000	Scale:	1.0000	
No. iterations:	6.0000			
	Coeff.	Std. Err.	z	P> z  [0.025 0.975]
const	-3.1110	1.0023	-3.1037	0.0019 -5.0755 -1.1464
age	0.0614	0.0144	4.2504	0.0000 0.0331 0.0896
ejection_fraction	-0.0734	0.0171	-4.3032	0.0000 -0.1069 -0.0399
serum_creatinine	0.7433	0.1756	4.2331	0.0000 0.3991 1.0874

```
In [441]: # Predicting probabilities on Test Data
y_pred_df = pd.DataFrame( { "actual": y_test,
    "predicted_prob": final_logit.predict(sm.add_constant( X_test[significant_
vars] ) )
    })
y_pred_df.predicted_prob.head(8)
```

Out[441]:

281	0	0.563221
265	0	0.125246
164	0	0.149884
9	0	0.998004
77	0	0.070533
238	0	0.151192
93	0	0.496610
109	0	0.124214

```
In [442]: from sklearn import metrics

# Roc Curve
```

```
# Roc Curve can be used to understand the overall performance of a logistic regression model and used
for model selection
# Roc Curve is a plot between False positive rate [FP / (TN + FP)] and True positive rate [TP / (TP+
FN)]
# Higher the AUC, better the model
# AUC < 0.5 model is inferior to a case of having no model
# AUC > 0.7 model is useful

def draw_roc( actual, probs ) :
    fpr, tpr, \
    thresholds = metrics.roc_curve( actual, probs, drop_intermediate = False )
    auc_score = metrics.roc_auc_score( actual, probs )
    plt.figure(figsize=(8, 6))
    plt.plot( fpr, tpr, label='ROC curve (area = %0.2f)' % auc_score )
    plt.plot([0, 1], [0, 1], 'k--')
    plt.xlim([0.0, 1.0])
    plt.ylim([0.0, 1.0])
    plt.xlabel('False Positive Rate or [1 - True Negative Rate]')
    plt.ylabel('True Positive Rate')
    plt.legend(loc='lower right')
    plt.show()
    return fpr, thresholds
```

```
fpr, tpr, thresholds = draw_roc( y_pred_df.actual, y_pred_df.predicted_prob)
```

```
auc_score = metrics.roc_auc_score( y_pred_df.actual, y_pred_df.predicted_prob )
print("auc_score", round( float( auc_score ), 2 ) )
```

True Positive Rate

False Positive Rate or [1 - True Negative Rate]

ROC curve (area = 0.71)

auc\_score: 0.71

```
In [443]: # Identifying the right probability
plt.figure( figsize = (8,6) )
sns.kdeplot( y_pred_df[y_pred_df.actual == 1]['predicted_prob'], color = '#31688e', shade=True, label =
'Patient Died' )
sns.kdeplot( y_pred_df[y_pred_df.actual == 0]['predicted_prob'], color = '#440154', shade=True, label =
'Patient Lived' )
plt.legend()
plt.show()
```

True Positive Rate

False Positive Rate or [1 - True Negative Rate]

ROC curve (area = 0.71)

auc\_score: 0.71

```
In [444]: y_pred_df["predicted"] = y_pred_df.predicted_prob.map(lambda x: 1 if x > 0.25 else 0)
print("Y Predicted")
print(y_pred_df.predicted.value_counts())
print("")
print("Y Actual")
print(y_pred_df.actual.value_counts())
print("")
print("Confusion matrix")
from sklearn.metrics import confusion_matrix
cm = metrics.confusion_matrix( y_pred_df.actual, y_pred_df.predicted, [0,1] )
print(cm)
```

labels = ["TN","FP","FN","TP"]

categories = ["Patient Lived [0]", "Patient Died [1]"]

make\_confusion\_matrix(cm,

group\_names=labels,

categories=categories,

cmmap="viridis\_r", figsize=(8,6))

Y Predicted

0 34

1 26

Name: predicted, dtype: int64

Y Actual

0 35

1 25

Name: actual, dtype: int64

confusion matrix

[[24 11]

[10 15]]

TN 24 40.00%

FP 11 18.33%

FN 10 16.67%

TP 15 25.00%

Accuracy=0.650

Precision=0.577

Recall=0.600

F1 Score=0.588

```
In [450]: y_pred_df["predicted"] = y_pred_df.predicted_prob.map(lambda x: 1 if x > 0.23 else 0)
print("Y Predicted")
print(y_pred_df.predicted.value_counts())
print("")
print("Y Actual")
print(y_pred_df.actual.value_counts())
print("")
print("Confusion matrix")
from sklearn.metrics import confusion_matrix
cm = metrics.confusion_matrix( y_pred_df.actual, y_pred_df.predicted, [0,1] )
print(cm)
```

labels = ["TN","FP","FN","TP"]

categories = ["Patient Lived [0]", "Patient Died [1]"]

make\_confusion\_matrix(cm,

group\_names=labels,

categories=categories,

cmmap="viridis\_r", figsize=(8,6))

Y Predicted

0 40

1 20

Name: gnb\_predict, dtype: int64

Y Actual

0 35

1 25

Name: actual, dtype: int64

confusion matrix

[[29 6]

[11 14]]

TN 29 48.33%

FP 6 10.00%

FN 11 18.33%

TP 14 23.33%

Accuracy=0.717

Precision=0.700

Recall=0.560

F1 Score=0.622

```
In [522]: y_pred_df["gnb_predict"] = y_pred_df.predicted_prob.map(lambda x: 1 if x > 0.18 else 0)
print("Y Predicted")
print(y_pred_df.gnb_predict.value_counts())
print("")
print("Y Actual")
print(y_pred_df.actual.value_counts())
print("")
print("Confusion matrix")
from sklearn.metrics import confusion_matrix
cm = metrics.confusion_matrix( y_pred_df.actual, y_pred_df.gnb_predict, [0,1] )
print(cm)
```

labels = ["TN","FP","FN","TP"]

categories = ["Patient Lived [0]", "Patient Died [1]"]

make\_confusion\_matrix(cm,

group\_names=labels,

categories=categories,

cmmap="viridis\_r", figsize=(8,6))

Y Predicted

0 40

1 20

Name: gnb\_predict, dtype: int64

Y Actual

0 35

1 25

Name: actual, dtype: int64

confusion matrix

[[29 6]

[11 14]]

TN 29 48.33%

FP 6 10.00%

FN 11 18.33%

TP 14 23.33%

Accuracy=0.717

Precision=0.700

Recall=0.560

F1 Score=0.622



```
[517]: y_pred_df['gnb_predict'] = y_pred_df.predicted_prob.map(lambda x: 1 if x > 0.20 else 0)

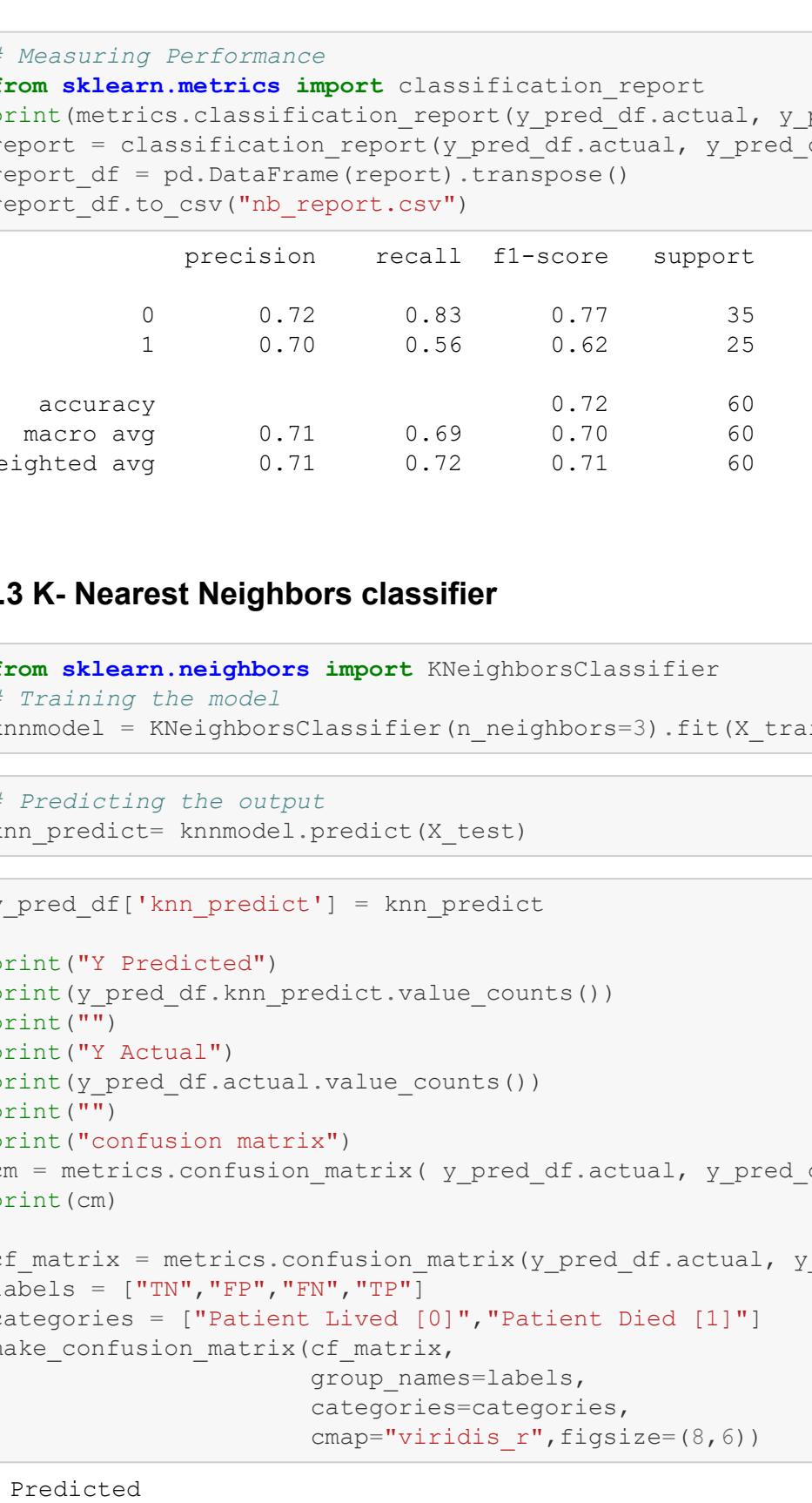
print("Y Predicted")
print(y_pred_df.gnb_predict.value_counts())
print("")
print("Y Actual")
print(y_pred_df.actual.value_counts())
print("")
print("Confusion matrix")
cm = metrics.confusion_matrix( y_pred_df.actual, y_pred_df.gnb_predict, (0,1) )
print(cm)

cf_matrix = metrics.confusion_matrix(y_pred_df.actual, y_pred_df.gnb_predict)
labels = ["TN", "FP", "FN", "TP"]
categories = ["Patient Lived (0)", "Patient Died (1)"]
make_confusion_matrix(cf_matrix,
                      group_names=labels,
                      categories=categories,
                      cmap="viridis_r", figsize=(8,6))

Y Predicted
0    42
1     18
Name: gnb_predict, dtype: int64

Y Actual
0     35
1     25
Name: actual, dtype: int64

Confusion matrix
[[30  5]
 [12 13]]
```



Accuracy=0.717  
Precision=0.722  
Recall=0.520  
F1 Score=0.605

```
In [523]: # Measuring Performance
from sklearn.metrics import classification_report
print(metrics.classification_report(y_pred_df.actual, y_pred_df.gnb_predict))
report = classification_report(y_pred_df.actual, y_pred_df.gnb_predict, output_dict=True)
report_df = pd.DataFrame(report).transpose()
report_df.to_csv("nb_report.csv")
```

	precision	recall	f1-score	support
0	0.72	0.83	0.77	35
1	0.70	0.56	0.62	25
accuracy			0.72	60
macro avg	0.71	0.69	0.70	60
weighted avg	0.71	0.72	0.71	60

### 3.3 K- Nearest Neighbors classifier

```
In [497]: from sklearn.neighbors import KNeighborsClassifier
# Training the model
knnmodel = KNeighborsClassifier(n_neighbors=3).fit(X_train, y_train)
```

```
In [498]: # Predicting the output
knn_predict= knnmodel.predict(X_test)
```

```
In [502]: y_pred_df['knn_predict'] = knn_predict

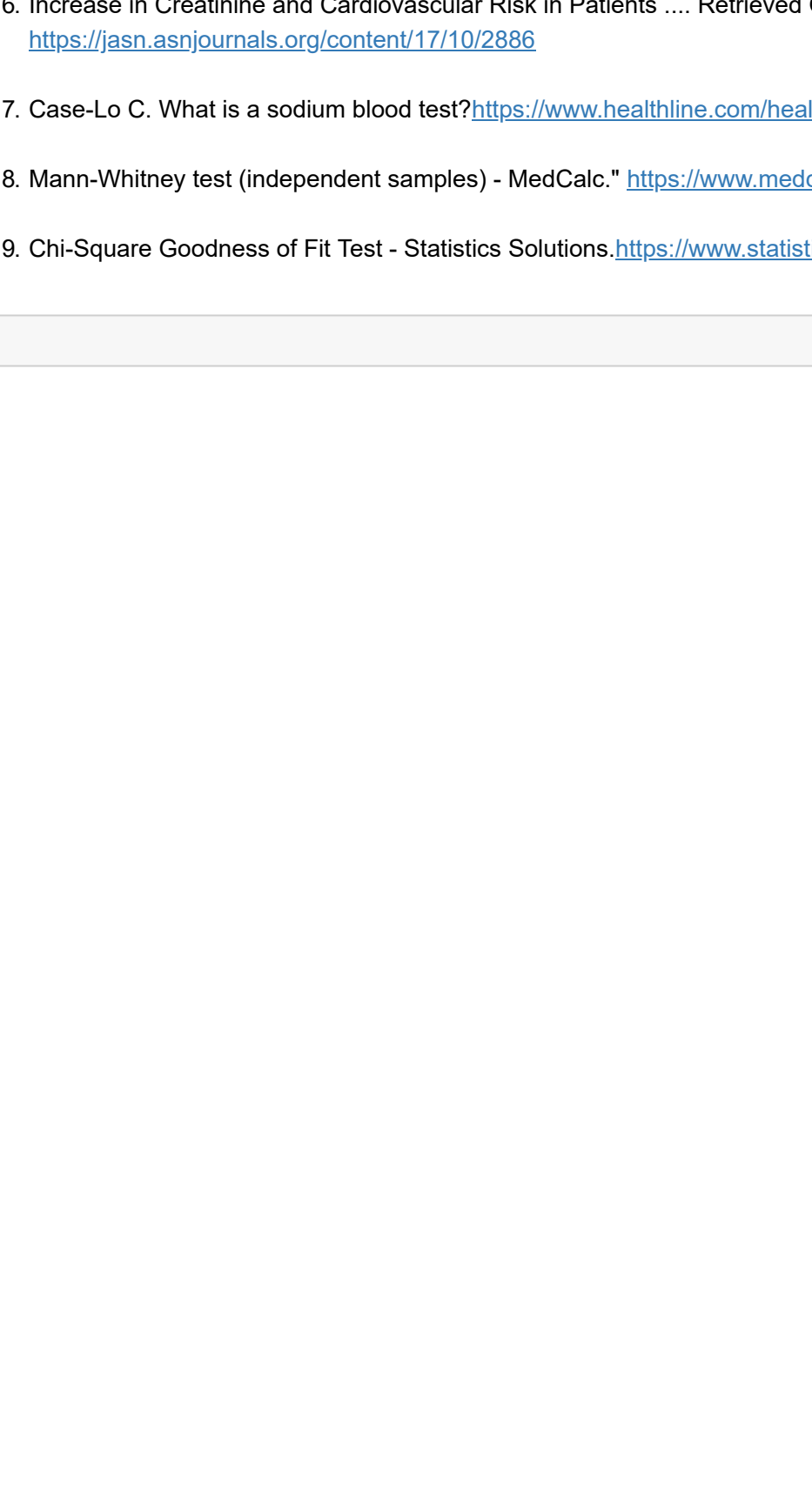
print("Y Predicted")
print(y_pred_df.knn_predict.value_counts())
print("")
print("Y Actual")
print(y_pred_df.actual.value_counts())
print("")
print("Confusion matrix")
cm = metrics.confusion_matrix( y_pred_df.actual, y_pred_df.knn_predict, (0,1) )
print(cm)

cf_matrix = metrics.confusion_matrix(y_pred_df.actual, y_pred_df.knn_predict)
labels = ["TN", "FP", "FN", "TP"]
categories = ["Patient Lived (0)", "Patient Died (1)"]
make_confusion_matrix(cf_matrix,
                      group_names=labels,
                      categories=categories,
                      cmap="viridis_r", figsize=(8,6))

Y Predicted
0     49
1     11
Name: knn_predict, dtype: int64

Y Actual
0     35
1     25
Name: actual, dtype: int64

Confusion matrix
[[28  7]
 [21 11]]
```



Accuracy=0.533  
Precision=0.364  
Recall=0.160  
F1 Score=0.222

```
In [503]: print( metrics.classification_report( y_pred_df['actual'], y_pred_df['knn_predict'] ) )
```

	precision	recall	f1-score	support
0	0.57	0.80	0.67	35
1	0.36	0.16	0.22	25
accuracy			0.53	60
macro avg	0.47	0.48	0.44	60
weighted avg	0.48	0.53	0.48	60

### 3.4 Results

**Logistic Regression:** AUC - 0.71 ; Accuracy 0.68 ; Sensitivity 0.68 ; Specificity 0.69 ; Precision 0.61  
**N-Bayes :** AUC - 0.73 ; Accuracy 0.72 ; Sensitivity 0.56 ; Specificity 0.83 ; Precision 0.72  
**KNN :** AUC - NA ; Accuracy 0.53 ; Sensitivity 0.16 ; Specificity 0.80 ; Precision 0.36\*\*

### References

- Chicco, D., Jurman, G. Machine learning can predict survival of patients with heart failure from serum creatinine and ejection fraction alone. BMC Med Inform Decis Mak 20, 16 (2020). <https://doi.org/10.1186/s12911-020-1023-5>
- Ahmad T, Munir A, Bhatti SH, Aftab M, Raza MA. Survival analysis of heart failure patients: a case study. PLoS ONE. 2017; 12(7):0181001.
- Johns Hopkins Rheumatology. Creatine Phosphokinase (CPK). <https://www.hopkinslupus.org/lupus-tests/clinical-tests/creatinine-phosphokinase-ckp/>. Accessed 25 Jan 2019.
- Ejection Fraction Heart Failure Measurement | American Heart <https://www.heart.org/en/health-topics/heart-failure/diagnosing-heart-failure/ejection-fraction-heart-failure-measurement>
- Platelets and Cardiovascular Disease - AHA Journals <https://www.ahajournals.org/doi/full/10.1161/01.CIR.0000086897.15588.4B>
- Increase in Creatinine and Cardiovascular Risk in Patients ... Retrieved October 10, 2020, from <https://aasn.sagepub.com/content/17/10/2698>
- Case-Lo C. What is a sodium blood test? <https://www.healthline.com/health/sodium-blood>.
- Mann-Whitney test (independent samples) - MedCalc. <https://www.medcalc.org/manual/mannwhitney.php>.
- Chi-Square Goodness of Fit Test - Statistics Solutions <https://www.statisticssolutions.com/chi-square-goodness-of-fit-test/>

```
In [ ] : 
```